# Theinteractionbetweenlocalfocusingstructureandglobalintentionsin spokendiscourse

SofiaGustafson -Capková
DepartmentofLinguistics,ComputationalLinguistics
StockholmUniversity
S-10961,Stockholm
Sweden
sofia@ling.su.se

ABSTRACT

Thepurposeofthestudyreportedinthispaperistoinvestigatehowlocalfocusingstructure,analysedintermsofCenteringTheory(Grosz, Joshi&Weinstein,1995),andglobald       iscoursestructure,analysedintermsofdiscoursesegmentsanddiscoursesegmentpurposes(Grosz& Sidner,1986),interact.SwedishdialoguewasanalysedaccordingtoCenteringTheoryandGroszandSidners(1986)discoursetheory.The resultsindicatean   interactionbetweenlocallyimplicitelementsandglobalintentions.Alsoindicationsconcerningdiscoursemarkersvarying intonationwerefound.

## Introduction

Discoursecanbedescribedasbuiltupfromdiscoursebuildingblockscalleddiscoursesegments           (hereafterDS). TheseDSaretheunitsforformingahierarchicaldiscoursestructure.Theyaredescribedine.g.Grosz&Sidner (1986,hereafterG&S)whereclaimsaremadeabouttheuseoftheDSintheglobaldiscoursestructureaswellas theirconnect iontothecoherenceofthediscourse,andinCenteringTheory(Grosz,Joshi&Weinstein,1995, hereafterCT),whereclaimsaremadeabouttheinternalstructureandcoherenceoftheDS:s. Grosz&Sidner(1986)haveappliedtheirdiscoursetheorytoboth            argumentativetextandtaskorienteddialogue, whileCTtraditionallyhasbeenappliedtonarrativetext.Inrecenttimes,however,aninterestforapplyingCTto dialoguehasarisen,andsomeattemptstodothathasbeencarriedout(e.g.Brennan,1998,B           yron&Stent,1998, Eckert&Strube,1999).

InapplyingG&Stheory,aproblematicpointistheimportanceofspeakerintention,whichgovernsboththe discoursesegmentingandthediscoursestructure.Itisunclearwhoseperspective,thatshouldbetaken;         thespeakers originalintention,thelistenersunderstandingoftheintentionortheanalysersinterpretationoftheintention. However,onethingthatisforsureisthattheanalyserwillcertainlyfaceachallengeifattemptingtofindout originalspe aker/listenerintentions.

ThemajorproblematicissuesinapplyingCTtheissueofbothutterancesegmentinganddiscoursesegmenting effectsalmostallotheraspectsoftheanalysis.AnotherproblemwithCTistodecidewhatconceptsthatare accessible,or *realised*,inanutterance.

Thepilotstudypresentedinthispaperisofexplorativecharacter,andaddressesarangeofproblemsencounteredin acombinedG&S -typeanalysisandCTanalysisoftaskorienteddialogue,i.e.:

- Utterancesegmenting,i.e .theunitsbetweenwhichlocalcoherenceiscomputed
- Discoursesegmenting,i.e.largerconstituentsaffectingtheglobaldiscoursestructure.Thesesegments correlateswithwhatCarlettaetal.(1997)calls"game".
- Whatitemsthatarepossiblecent  ers
- TheCTnotionofa *realised*item

Theaimofthispaperistogiveapictureofhowthoseproblemsareconnectedtoeachother,andtooutlinehowto refineamultiple -levelanalysis.Itisalsoanattempttoapplyaglobalandlocalanalysistospok          enlanguagedata, andtogiveaccountforspecificproblemsthatarisesbysuchananalysis.Itisthehopeoftheauthorthatresults frominvestigationslikethisshouldhelptodevelope.g.instructionsformoreextensiveinvestigationsinthefield.

## 1   Background

G&SandCTaretwotheoriesthatgiveaccountfordiscoursestructureandcoherence,buta            tdifferentlevelsofthe discourse.G&Smainlyaddressestheglobaldiscoursestructure,whileCTgivesaccountforthelocalcoherence.I willheregiveashortdescriptionofboththeories.

G&S(1986)describediscourseasconsistingofthreestructures           :i)thelinguisticstructure,ii)theintentional structureandiii)theattentionalstate.Thesethreestructuresinteract,buttheyarestilltobeconsideredasseparate structures.Theinteractionbetweenthemworksroughlyasfollow:Thelinguistics          tructure,i.e.thestringofwords, isdividedintodiscoursesegments.EachsegmenthasaDiscourseSegmentPurpose(DSP)whichispartofthe intentionalstructure.AccordingtohowtheDSP:saresatisfied,differentrelationsholdbetweenthediscourse segmentsandtheattentionalstateismodelledoutoftheserelations.

Thus,inG&S,disoursesegmentsareintentionallydelimited,i.e.adiscoursesegmentisgovernedbyamain intention,theDSP.TherangeofDSP:sisunlimited.DSmaybenested,an            dtherelationsthatholdbetween discoursesegmentsarelimitedtotwo:i)dominanceandii)satisfaction -precedence.Dominancemeansinshortthat adiscoursesegmentBwhichispartofthesatisfactionoftheintentiongoverningthediscoursesegmentA              is dominatedbyA,i.e.AdominatesB.Satisfaction         -precedenceontheotherhandholdsinthecaseswherethe

intentionofadiscoursesegmentChastobefulfilledbeforetheintentionofthediscoursesegmentDappears,i.e.C satisfaction-precedesD.

Therelationsdominanceandsatisfactionprecedencecontributesincrementallytothediscoursestructureandmodel theglobalcoherence.Thisisdonebystackmanipulations,whichcouldbedescribedasmodellingatemporal sequenceofintentionsinattention alfocusinthediscourse.Thisprocesswillhowevernotbecloselydescribed here.

CenteringTheoryisatheory,whichgivesaccountforthedegreeoflocalcoherencebetweenutteranceswithina discoursesegment.Thisismadebysegmentingthelinguistic stringintoutterancesandclassifythetransitions betweenthem.Thetransitionsarecomputedonbasisoftwofactors:backward -lookingcenterCbandforward -lookingcenter(s)Cf.Sometimesthepreferredcenter,i.e.thehighestrankedmemberoftheCf -listissingledoutas Cp.ThechoiceofcentersisstandardCT(Grosz,Joshi&Weinstein)basedongrammaticalroles: subject>object>otherroles.Itisimportanttonote,thatthecentersdoesnothavetobeexplicitlypresentinthe linguisticstring( directlyrealized),butmayalsobeimplicitlypresentintheconceptualrepresentation(realized). Thismeans,thatcentersarenotlinguisticunits,butconcepts.

ThefourtransitionsarecomputedonbasisoftheC:s,asshownin Table 1.

**Table 1TableoverthetransitionsinCenteringTheory.**

|  | $Cb(U_i)=Cb(U_{i-1})$ or $Cb(U_{i-1})=[?]$ | $Cb(U_i)\neq Cb(U_{i-1})$ |
|---|---|---|
| $Cb(U_i)=Cp(U_i)$ | CONTINUE | SMOOTH-SHIFT |
| $Cb(U_i)\neq Cp(U_i)$ | RETAIN | ROUGH-SHIFT |

InadditiontworulesareusedinCT.Thefirstis"Thepronounconstraint".ThisrulestatethatifsomethinginU $_i$is realizedasapronouninU $_{i+1}$theCbofU $_{i+1}$mustalsoberealisedwithapronoun.

Thesecondrulestatesthatsequencesofcontinuation arepreferredoversequencesofretain.Theshifttransitions putgenerallyahigerinferenceloaduponthehearer.

ThesegmentingissueiscertainlyimportantalsoinCT,butitisnotcloseraddressedbyGrosz,Joshi&Weinstein (1995),i.e.noexplicit baseforthediscoursesegmentsisgivenhere.Itishoweveragoodguessthattheyshouldbe ofthesamenatureasbyGrosz&Sidner,whomentionscenteringaspossibleadditionalmechanism(Grosz& Sidner,1986,p.91).

AmodifiedversionofCTismade byWalker(1997).ShehasreplacedtheDSbysomethingthatcouldbedescribed asamovingwindow.ThismeansthattheCT -analysisisdonecontinuallythroughthewholediscourse,anditdoes notstartandstopoverandoveragainbytheinitiationorendi ngofaDS.TheelementsfromtheCf -listaresavedin acache,whichisincrementallyupdatedinthewaythatnewitemsareaddedandolditemsareerased.Walker suggeststhatthesizeofthemovingwindowshouldconsistoftwoorthreesentences,orsev enpropositions. Bothutteranceboundariesanddiscoursesegmentboundariesaredifficulttodelimitinspokenlanguage.Utterances aredifficultbecausethereisoftennoformallycorrectlycompletedsentencestructureinspontaneousspeech. Generalstrat egiesforsegmentationofspokendiscourseareprosodicphrasing,cue -words,andtheuseofformfor referringexpressions(e.g.Passonneau&Litman,1997,Grosz&Sidner,1986,Walker,1997).

Anadditionalprobleminanalysingdialogueisthatitisno tquiteclearhowtoapplyatheorylikeCT,mainly developedwithworkonnarrativetext,foramulti -partydiscourse.E.g.isthepreviousutteranceforXthelinearly previousutterance,orthepreviousutteranceutteredbyX?Onehastoworkwithatle asttwopersonsinterpretations ofthediscourse,interpretationswhichdonothavetobeoverlapping,intermsofbothDSandfocusofattention, i.e.oneshouldtrytokeeptrackonwhosecenterthatisanalysed.

## 2    Method

InordertogivespokenlanguageadiscourseanalysisintermsofbothG&SandCT,elicitedspokendialoguewas analysed.ThespokenlanguagematerialwasfromoneMapTaskdialogueinSwedish(Helgason).Inall60turns fromonedialoguewithtwospeakers wereanalysed.

Thedialoguewassegmentedwithapause -detectingtool,whichdetectedsilentpauseslongerthan100ms.After examinationofthesegmenteddatatheanalyserdecidedthatpauses300msorlongershouldbeusedasutterance boundaries.This pauselengthisroughlycorrelatingwithclauseboundariesaccordingtoGarman(1990),whosets clauseboundariesto400ms.Thesegmentationbasedon300msorlongerpausesresultedin100utterances.The decisionwasalsomadethatchangeofspeakerals oindicatednewutterance.

Alllinguisticunitswereregardedasvalid,i.e.nofilteringoutofutterancesconsistingonlyofe.g.humming (mmm…)wasdone,asdonebye.g.Byron&Stent(1998).

Thetransitionsbetweentheutteranceswerecomputedaccordi ngtoCT,buttherankingofcenterswerelimitedto linearappearance.Whenitcomestoitemspossiblecarrycenter,1and2sg.pronounswerefilteredout.Afterthis thematerialwasexaminedaccordingtointentionalcontent.BoundariesbetweenDS,whic hcorrelatedtocertain intentions,wereannotatedandalsotheintentionsweredescribed.Thisresultedin36labelleddiscoursesegments, whichwereanalysedintermsofrelationsbetweendifferentDS(dominanceandsatisfaction -precedence).Change ofs peakerwasnottakentoimplynewDS.

2

# 3 Resultsanddiscussion

Oneofthemainproblemsindiscourseanalysisischoosinganinterpretationthatisasgeneralaspossible,i.e.totry tominimisethesubjectivityintheinterpreta tion.Thereasonforthisistheneedforapossiblereplicationofthe analysisoftheinterpretation,i.e.theanalysisshouldnotbetooboundedtotheanalyserssubjectivelybased interpretation.Theanalyseristhusforcedtokeeplanguageinterpretat ionanddiscourseanalysisstrictlyseparated. Thisisinitselfaparadox,becausetoanalyseastretchofdiscoursemeanstoanalyseaninterpretationofthestretch ofdiscourse.However,humannaturallanguageisneverimpersonallyinterpreted,itisa lwaysinterpretedthrough thefilterofasubjectivehumanthinking,sotheinterpretationandtheanalysisblend.Infieldsase.g.computational linguistics,onetriestomodelapureandobjectiveinterpretation,whichinfactisthemostunlikelyinter pretationin itspureness.Thequestionis,howisitpossibletokeeponsearchingforthemostgeneralinterpretationbutstill avoidbothsubjectivityandartificialityintheinterpretation,i.e.tomaketheclaimoftheanalysispartofthe interpretation,asobjectiveaspossible,sothatitisscientificallyvalid,butstillkeepasmuchsubjectivityas possibleintheinterpretationsothattheoutcomemimicslanguageusersasmuchaspossible.Inthesubjectivityof humanlanguageunderstandinglie salsotherobustnessandthegeneralityofhumanlanguageuse. Tousedialogueisonewaytotrytodelimitthedegreeofsubjectivityintheanalysis,butstillallowsubjectivityin theinterpreta tion.Thereasonforthisisthattheprimarytaskfortheanalyserisnottointerpretthetext/speech,but tounderstandhowthecurrentspeakerinterpretedwhattheformerspeakersaid.Thismeansthattheanalyserhasa referencepointfortheinterpre tationoutsideherself.Infollowingthedialogueitisalsopossibletofollowhowa personactuallyinterpretsthecurrentspeaker.Theanalyserisnotcompletelyalonewithherowninterpretation,but isabletogetaglimpseofhowanotherpersoninter pretstheutterances.

Tousethepause -toolfordetectingutteranceboundarieswasanotherwaytotrytolimittheinfluenceofsubjective interpretation.Theinterpreterwasnotdeterminingthesegmentationherself,butusedakindofbootstrappingin decidingtheutteranceunits.

Theresultsshowedthattheutterancesegmentationinmanycaseswasquitegood,butstill,inmanycasesthe granularitywasfinerthanpreferred.Thediscoursesegmentingonintentionalbasisdidnotposegreatproblems, butperhapsthatjustindicatesthereadinessbytheanalysertoassignexplicitintentionstocertainsegments.Below anoverviewofthesegmentedmaterialisshown.

- Turns:60
- Utterances:100
- Discoursesegments:36

## 3.1 Segmentingtheu tterances

Caseswherethepause -basedutterancesegmentationwasnotoptimalcouldsometimeshavebeenavoidedifthe intonationcontourhadbeentakenintoconsideration.In Example 1,givenbelow,the speechsignalwassegmented atthepointUtt2.(thepauseprecedingthatpositionislongerthan300ms),butthisbreakcouldhavebeenavoided ifthefact,thattheintonationcontourisstable(i.e.neitherrisingnorfalling)hadbeentakenintocons ideration.

**Example 1**

Utt.1..dåskavisedåharvi..en..ens!0.; *karta*härframföross..ochjaghar..;. *landstigit*påen *plats*, →
Utt.2..pådenhärön.

Intheana lysiscaseslike Example 1werehoweverregardedastwoutterances.

## 3.2 TheCTanalysis

Afterthesegmentationintoutterances,the100utteranceswereanalysedintermsofCT.Atthispoi ntinthe investigationnoattempttodividethediscourseintoDSwasmade.FollowingWalker(1997)acontinual examinationofthecentersandthetransitionswasdonethroughoutthewholediscourse.

### 3.2.1 Analysingthecenters

Concerningtheanalysisofthecenterstherankingbasedongrammaticalroledidnotturnouttobesuitableforthe analyseddialogue,partlyduetothefactthat1and2sg.pronounswerefilteredout.Insteadtheanalyserfollowed threesimplestatements :

- Allkindsofelements(e.g.complexphrasesaswellassinglewords)wererankedafterthelinearoccurrencein thespeechsignal.
- Phoneticprominencewastakenintoaccount.Aphoneticallyprominentelementwasgivenahigherrankthana phoneticallynon -prominentelement(thepronounconstraintwashoweveralwayskept).
- Coordinateandsubordinateclauseswerespeciallyhandled.Acon -/subjunctioninsideanutterancestarteda newCf -list,whichwasgivenhigherprominencethanthefirstlist.

Anexampleoftherankingofelementsinanutteranceisgivenin
Example 2,wheretheunderlined"men"initiatesthesecondCf -list(alsounderlined),fromwherethepreferred centerischosen.

**Example 2**

.ja..detärett ***aningers..närnär*** mareflodenän;..;..kust!0..!0kantendär     <u>men</u>detär ***nästanmitt*** emellan.

Cb=ja<dusnuddarnästanvidenflodnärduärdär(precedingutterance)>
Cf=1.[närmare  floden,kustkanten]2.  <u>[mittemellan<floden&kustkanten>]</u>
Cp=mittemellan<floden&kustkanten>


Asearliernoted,Cb:scouldbedirectlyrealisedorrealised.In          Example 1forinstancesomeelementsarepresentin theanalysis,butnotpresentintheutteranceorintheappropriateplaceintheutterance.Suchinstancesaremarked outwith<>intheanalysis.In        Example 1therearet woinstancesofsuchpartlyimplicitelements:1.<dusnuddar nästanvidenflodnärduärdär>and2.<floden&kustkanten>.Inthefirstcase,"ja"doesnotonlyseemtobea waytosignalthatthelistenerhaveunderstand,butalsoawaytosignalthat           therepresentationoftheconceptsis stillrelevantandactive,i.e."ja"functionsasashort"keepactive         -signal".Thisisfoundinallja -instances.Similar findingsarereportedbyEckert&Strube(1999),whoclaimthatthoseutteranceshavehighrel       evanceforgrounding indialogue.Inthesecondcaseboththeriver(flod)andtheshore(kustkanten)areintroduced,butinthelaterpartof theutterance(after"men")thefocusisonthepointbetweenthebothelements.However,theelementsarestill highlyactiveindefiningthepointinbetweenthatiswhytheconjunctionoftheboothconceptsisanalysedas present.Suchpartlyimplicitelementsarefrequentinthematerial,butalsocompletelyimplicitelements. Interestingis,thatinthecaseof      acompleteimplicitelementinanutteranceas:".och      ***fortsätter***<vägen>norrut", theCpintheutterance,theconcept"road"(<vägen>)isacrucialconceptintheformulationoftheDSP(the discoursesegmentpurposes,theintentionsmotivatingadiscours       esegment).Thus,theconceptcouldbesaidtobe contextuallyhighlyactivated,i.e.activatedbythetaskandthesituationitself,oractivatedonthegloballevel.As wellaswecantalkaboutlocalandglobalfocus,wecanalsobeabletodistinguish         betweenalocalandaglobal levelofactivation.
Togetaviewovertheproportionsofimplicitvs.explicitreferenceindiscoursealltheCb:swerecountedand sortedasdirectlyrealised(explicitlypresent)orrealised(implicitorpartlyimplicitpre        sent).Theresultisshownin Table 2.

**Table 2**

|        | Explicitlypresent | partlyimplicitlypresent | fullyimplicitlypresent |
|--------|-------------------|-------------------------|------------------------|
| N=100  | 19                | 71                      | 10                     |

Theabovefiguresindicatethat81%ofbackreferenceinadiscourseisimplicit,whichmakeshuman communicationseemlikeaniceberg.
Theproportionoftransitionsbetweenutteranceswascomputed,andtheresultsaregivenin       Table 3.

**Table 3Transitionsbetweenthe100utterancesinthematerial**

| Continue | Retain | Smooth-shift | Rough-shift |
|----------|--------|--------------|-------------|
| 47       | 36     | 10           | 5           |

PleasenotthatoneinstanceofRough      -shiftisclearlydiscourseinitial,soitisleftoutinthetableabove.These resultswillbecloserdiscussedundertheheading       3.3.

### *3.3    Theglobalstructureofthediscourse*

Theglobalstructureofthediscoursewasanalysedintermsoftherelationsdominanceandsatisfaction        -precedence betweendiscoursesegments.Inmakingthisanalysistheanalyserexperiencedaneedtomakeamorefine        -grained distinctionbetweendifferentinstancesoftherelationsatisfaction     -precedence.Thus,therelationsusedwere:
- Dominance,correspondstot hedominancerelationbetweentwosegments(mother     -daughter).
- Singlepop,correspondstotwoadjacentsegmentsonthesamelevelbothwithoutdaughters(sisterswithout daughters).Thisisthesameastherelationsatisfaction       -precedencebetweentwosi sterswithoutdaughters.
- Multiplepop,correspondstotwosegments,textuallyadjacentbutondifferentlevelsinthehierarchical analysis,i.e.theyoungestdaughterinonebranchandapotentialmotherforanotherbranch.Thisisthesameas satisfaction-precedencebetweentwonodesondifferenthierarchicallevels.

Asnotedin    Table 3abovetherewasanoverwhelmingnumberoftherelationsContinueandRetain.Bothshift        - transitionswerequiterare,anditis      worthnotingthatallRough      -shiftsappearedeitheri)insideadiscoursesegment (4)ii)betweentwodiscoursesegmentsrelatedtoeachotherbytherelationdominance(1)iii)afteraveryclear indicationthatthediscoursetopicwillchange("dåså,då       skavise".Thisistheoneleftoutin        Table 3).Thelast alternativeispossibletoexcludeonthebasisthatitisbettertoconsiderthisasanewdiscourseandnotashift insidethesamediscourse.Thetwofi rstalternativeshoweverindicate,thatrough      -shiftappearsonlyinsideatightly definedintentionalspace,inthedataitneverappearstogetherwithashiftoftheintention.Itneverappeared betweenDS:srelatedwithSinglepoporMultiplepop.Amech          anismasrough     -shiftseemsthusnottobethe

appropriatewaytomakesuchachangeofdirectioninthediscourse,ratheritindicatesmisunderstandingora
"jumping"insideoneisolatedintentionalspace.
Thetransitionsbetweendiscoursesegmentsissho    wnin Table 4below.

**Table 4CT -TransitionsatdifferentkindsofDSboundaries.**

|              | Continue | Retain | Smooth-shift | Rough-shift |
|--------------|----------|--------|--------------|-------------|
| Multiplepop  | 0        | 9      | 0            | 0           |
| Singlepop    | 0        | 6      | 2            | 1           |
| Dominance    | 6        | 6      | 0            | 4           |

Ininvestigatingwhatcouldbecharacteristicfordiscoursesegmentboundarieswithdifferentrelationsalldiscourse
segmentboundarieswereinvestigated.Theresultsaregivenbelow.

- Multiplepop:Indicatesinsevencasesofninewithacombinationofpause,cue                -wordsandphonetic
  prominence( *ochsen* /*ochfortsätter* ).
- Singlepop:Indicatesinfivecasesofninewithacombinationofpauseandcue          -words(ochsen/då)
- Dominance: Nospecialpreferencesfound.

# 4   Summaryandfurtherwork

Thisinvestigationreportedinthispaperiscertainlysufferingfromarangeofweakpoints;forinstancealargerset
ofdataandanevaluationofintercoderreliability          wouldbehighlydesirable.Theanalysisisnowverydependenton
oneanalysersowninterpretation.Theresultshowevergivequiteinterestingindicationsconcerningtheinteraction
betweenlocalfocusandglobalintentions,e.g.theconnectionbetweenthe          implicitcentersandtheintentionsbehind
thediscoursesegments.
Theuseofpausesforutterancesegmentationwouldcertainlybebetteriftheintonationcontourcouldbeintegrated
intheanalysis.Inthedataitalsoseemstobearegularityintheuse          ofintonationbytheuseofcue          -words,i.e.the
alternationbetweenphoneticallyprominentandphoneticallynon          -prominentcorrelateswiththedifferentrelations
MultiplepopandSinglepop,itishoweverdifficulttosayanythingforsurewithoutanalys          ingalargeramountof
data.
Furtherworkinthisdirectionwould,exceptmoredata,includeamorethoroughinvestigationoftherankingorder.
Toisolatewhatconceptsthatarepresent,orratheraccessible,inanutteranceinacertaincontextisalso          indeedan
important,butdifficulttasktoattack.Itwouldalsobeofinteresttoconnectfindingsfromanalysislikethisto
dialoguecoding,asdescribedbye.g.Carlettaetal.

# 5   Acknowledgments

# 6   Literature

**Carlettaetal.(1997):** *TheReliabilityofaDialogueCodingStructureScheme* .ComputationalLinguistics,vol.23
no.1.
**Berthelsen,H** .:Pausedetectingtool.
**Brennan,S(1998)** **:** *Centeringasapsychologicalresourceforachievingjointreferenceinspontaneousdiscourse* .
InWalkeretal.1998,CenteringTheoryindiscourse.
**Byron,D .&Stent,A.(1998):** *APreliminaryModelofCenteringinDialog.* IntheProceedingsofthe36thAnnual
MeetingoftheAssociationforComputationalLinguistics(ACL'98)studentsession.
**Eckert,M.&Strube,M(1999):** *ResolvingDiscourseDeicticAnaphorainDialogues* .InEACL '99.
**Garman,M.(1990):** *Psycholinguistics.* CambridgeUniversityPress.
**Grosz,B.&Sidner,C.(1986):** *Attention,IntentionsandtheStructureofDiscourse* .ComputationalLinguistics,
vol12,no3.
**Grosz,B.,Joshi,A.&Weinstein,S.(1995):** *Centering:AF rameworkforModellingtheLocalCoherenceof
Discourse* .ComputationalLinguistics,vol.21,no.2.
**Helgason,P.:** *TheStockholmCorpusofSpontaneousSwedish.* MapTaskcorpus,Departmentoflinguistics,
StockholmUniversity.
**Passonneau,R.&Litman,D. (1997):** *DiscourseSegmentationbyHumanandAutomatedMeans* .Computational
Linguistics,vol23,no1.